

2-12-02 #3

LAW OFFICES
SUGHRUE, MION, ZINN, MACPEAK & SEAS, PLLC
2100 PENNSYLVANIA AVENUE, N.W.
WASHINGTON, DC 20037-3213
TELEPHONE (202) 293-7060
FACSIMILE (202) 293-7860
www.sughrue.com

J. Frank Osha, Esq.
Direct Dial (202) 663-7915
Email: fosha@suehrue.com

March 14, 2001

BOX PATENT APPLICATION
Commissioner for Patents
Washington, D.C. 20231

Re: Application of Kenji YAMANISHI and Hang LI
QUESTIONNAIRE ANALYSIS SYSTEM
Our Ref. Q63084

Dear Sir:

Attached hereto is the application identified above including 30 sheets of the specification, including the claims and abstract, 10 sheets of drawings, the executed Assignment and PTO 1595 form, and the executed Declaration and Power of Attorney. Also enclosed is an Information Disclosure Statement with form PTO-1449 and references.

The Government filing fee is calculated as follows:

Total claims	13	-	20	=		x	\$18.00	=	
Independent claims	4	-	3	=	1	x	\$80.00	=	\$80.00
Base Fee									\$710.00
TOTAL FILING FEE									\$790.00
Recordation of Assignment									\$40.00
TOTAL FEE									\$830.00

Checks for the statutory filing fee of \$790.00 and Assignment recordation fee of \$40.00 are attached. You are also directed and authorized to charge or credit any difference or overpayment to Deposit Account No. 19-4880. The Commissioner is hereby authorized to charge any fees under 37 C.F.R. §§ 1.16 and 1.17 and any petitions for extension of time under 37 C.F.R. § 1.136 which may be required during the entire pendency of the application to Deposit Account No. 19-4880. A duplicate copy of this transmittal letter is attached.

Priority is claimed from March 15, 2000 based on Japanese Application No. 2000-071657. The priority document is enclosed herewith.

Respectfully submitted,
SUGHRUE, MION, ZINN,
MACPEAK & SEAS, PLLC
Attorneys for Applicant

By: J. Frank Osha
J. Frank Osha
Registration No. 24,625

日 本 国 特 許 庁

PATENT OFFICE
JAPANESE GOVERNMENT

Yamanishi et al
Filed 3/14/01

Q63084

10f1

J1017 U.S. PTO

09/805114



別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日

Date of Application:

2000年 3月15日

出 願 番 号

Application Number:

特願2000-071657

出 願 人

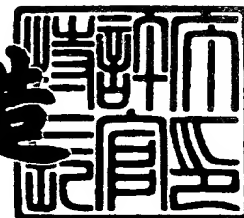
Applicant(s):

日本電気株式会社

2000年12月 1日

特許庁長官
Commissioner,
Patent Office

及 川 耕 造



出証番号 出証特2000-3100216

【書類名】 特許願

【整理番号】 33509710

【あて先】 特許庁長官殿

【国際特許分類】 G06F 17/00

【発明者】

 【住所又は居所】 東京都港区芝五丁目 7 番 1 号 日本電気株式会社内

 【氏名】 山西 健司

【発明者】

 【住所又は居所】 東京都港区芝五丁目 7 番 1 号 日本電気株式会社内

 【氏名】 李 航

【特許出願人】

 【識別番号】 000004237

 【氏名又は名称】 日本電気株式会社

【代理人】

 【識別番号】 100088890

 【弁理士】

 【氏名又は名称】 河原 純一

【手数料の表示】

 【予納台帳番号】 009690

 【納付金額】 21,000円

【提出物件の目録】

 【物件名】 明細書 1

 【物件名】 図面 1

 【物件名】 要約書 1

 【包括委任状番号】 9001717

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 アンケート回答分析システム

【特許請求の範囲】

【請求項 1】 自然言語による自由回答記述を含むアンケート回答文を入力する手段と、

前記アンケート回答文を送信するネットワークと、

前記送信されたアンケート回答文を蓄積するデータベースと、

該データベースから前記アンケート回答文を読み出して、前記アンケート回答文を分類するルールを学習するテキスト分類エンジンと

を備えることを特徴とするアンケート回答分析システム。

【請求項 2】 自然言語による自由回答記述を含むアンケート回答文を入力する手段と、

前記アンケート回答文を蓄積するデータベースと、

該データベースから前記アンケート回答文を読み出して、前記アンケート回答文を分類するルールを学習するテキスト分類エンジンと

を備えることを特徴とするアンケート回答分析システム。

【請求項 3】 自然言語による自由回答記述を含むアンケート回答文を入力する手段と、

前記アンケート回答文を送信するネットワークと、

前記送信されたアンケート回答文を蓄積するデータベースと、

該データベースから前記アンケート回答文を読み出して、前記アンケート回答文を分類するルールを学習するテキスト分類エンジンと、

該ルールを要求者からの要求に応じて前記ネットワークを通じて配信する手段と

を備えることを特徴とするアンケート回答分析システム。

【請求項 4】 前記テキスト分類エンジンが、前記データベースに蓄積されたアンケート回答文の全文に対して形態素解析を行う形態素解析手段と、カテゴリとテキストとを指定させるカテゴリ・テキスト指定手段と、前記データベースから読み込んだ複数のアンケート回答文に対して属性選択を行う属性選択手段と、前記

属性選択手段により属性選択された単語を基にテキストとカテゴリとの対応を表すルールを学習するルール学習手段と、前記ルール学習手段により学習されたルールを出力するルール出力手段とを含む請求項1、請求項2、または請求項3記載のアンケート回答分析システム。

【請求項5】前記属性選択手段が、テキスト中に出現する単語の各々について、該単語の出現を考慮しない場合のテキスト集合の確率的コンプレキシティと考慮する場合のテキスト集合の確率的コンプレキシティとの差 $\Delta SC(\omega)$ を計算する手順と、該差 $\Delta SC(\omega)$ がしきい値 τ より大きければ属性として選択する手順とを実行する請求項4記載のアンケート回答分析システム。

【請求項6】前記ルール学習手段が、テキストの集合を、 (d_1, c_1) , (d_2, c_2) , ..., (d_m, c_m) [ただし、各 d_i は多値の離散ベクトル $d_i = (\omega_{i1}, \omega_{i2}, \dots, \omega_{in})$ ($i = 1, \dots, m$)、 ω_{ij} は属性選択で得られた単語 ω_j ($j = 1, \dots, n$)が i 番目のテキストに現れれば1、そうでなければ0の値をとる、 c_i は i 番目のテキストに対応するカテゴリの値(ラベル)を表し、各 c_i は所与のカテゴリに属するならば1、そうでなければ0の値をとる、 m はテキストの数]の表現に置き換えて成形する手順と、拡張型確率的コンプレキシティ最小原理または確率的コンプレキシティ最小化原理などの情報量規準を用いてif-then-else型のルールを選択して確率的決定リストに順番に追加していく成長処理を行う手順と、前記確率的決定リストの最後からルールを1つずつ取り除き拡張型確率的コンプレキシティ最小原理の観点からみてこれを取り除かないほうがよいところまで続ける刈り込み処理を行う手順とを実行する請求項4記載のアンケート回答分析システム。

【請求項7】コンピュータを、前記データベースに蓄積されたアンケート回答文の全文に対して形態素解析を行う形態素解析手段、前記テキスト分類エンジンに対してカテゴリとテキストとを指定させるカテゴリ・テキスト指定手段、前記データベースから読み込んだ複数のアンケート回答文に対して属性選択を行う属性選択手段、前記属性選択手段により属性選択された単語を基にテキストとカテゴリとの対応を表すルールを学習するルール学習手段、および前記ルール学習手段により学習されたルールを出力するルール出力手段として機能させるためのプロ

グラムを記録した記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明はアンケート回答分析システムに関し、特にテキスト自動分類技術、自然言語処理技術、およびネットワーク利用技術を使用するアンケート回答分析システムに関する。

【0002】

【従来の技術】

インターネット等のネットワークを通じて得られる、自然言語による自由回答記述を含むアンケート回答文から全体的な特徴や傾向を導き出す作業は、従来の多くは人手で行われてきた。株式会社電通が開発したDE-FACTO(パンフレット)や富士通株式会社のKeyword Associator(渡部勇: 発散的思考支援システム「Keyword Associator」第二版, 計測自動制御学会第15回システム工学部会研究会資料, 1994年7月), HIPS等のテキストマイニングツール(渡辺, 三末, 新田, 杉山: ハイブリッド発想支援システム「HIPS」, 計測自動制御学会第15回システム工学部会研究会資料, 1995年1月)は、テキスト情報から言葉の関係性を導き出すのに用いられた。しかし、これらのツールは、ルールという形式でアンケート回答文の特徴を表現するものではなかった。

【0003】

また、自然言語による自由回答記述を含むアンケート回答文をインターネット等のネットワークを通じて自動的に収集して分析し、場合によっては要求者に分析結果を配信するようなシステムやサービスはこれまで存在しなかった。たとえば、特開平11-066036号公報, 特開平11-143856号公報等には、ネットワークを通してアンケートを実施し、回答をデータベースに蓄積する技術が開示されているが、ルールという形式でアンケート回答文の特徴を抽出することまでは行われていない。

【0004】

【発明が解決しようとする課題】

上述した従来の人手によるアンケート回答分析では、アンケート回答文の数が大量になると、非効率で手に負えないという問題点があった。

【0005】

また、DE-FACTOやHIPS等のテキストマイニングツールでは、ルールという形式でアンケート回答文の特徴を抽出できないので、コンパクトで分かりやすい知識の提示という点で十分ではないという問題点があった。

【0006】

さらに、検索等で用いられているテキスト分類ツールは、アンケート回答分析にも有力なツールであるが、自然言語による自由回答記述を含むアンケート回答文の分析に用いられている前例はなかった。

【0007】

本発明の第1の目的は、テキスト分類エンジンを用いることによって、自然言語による自由回答記述を含むアンケート回答文からコンパクトで分かりやすいルール形式の知識を自動的に提示することを可能にするアンケート回答分析システムを提供することにある。

【0008】

また、本発明の第2の目的は、上記に加えて、ネットワークを通じて集めたアンケート回答文からルール形式の知識を自動的に抽出し、要求者に分析結果を配信するアンケート回答分析システムを提供することにある。

【0009】

【課題を解決するための手段】

本発明のアンケート回答分析システムは、自然言語による自由回答記述を含むアンケート回答文を入力する手段と、前記アンケート回答文を送信するネットワークと、前記送信されたアンケート回答文を蓄積するデータベースと、該データベースから前記アンケート回答文を読み出して、前記アンケート回答文进行分类するルールを学習するテキスト分類エンジンとを備えることを特徴とする。

【0010】

また、本発明のアンケート回答分析システムは、自然言語による自由回答記述を

含むアンケート回答文を入力する手段と、前記アンケート回答文を蓄積するデータベースと、該データベースから前記アンケート回答文を読み出して、前記アンケート回答文を分類するルールを学習するテキスト分類エンジンとを備えることを特徴とする。

【0011】

さらに、本発明のアンケート回答分析システムは、自然言語による自由回答記述を含むアンケート回答文を入力する手段と、前記アンケート回答文を送信するネットワークと、前記送信されたアンケート回答文を蓄積するデータベースと、該データベースから前記アンケート回答文を読み出して、前記アンケート回答文を分類するルールを学習するテキスト分類エンジンと、該ルールを要求者からの要求に応じて前記ネットワークを通じて配信する手段とを備えることを特徴とする。

【0012】

【発明の実施の形態】

以下、本発明の実施の形態について図面を参照して詳細に説明する。

【0013】

(1) 第1の実施の形態

図1は、本発明の第1の実施の形態に係るアンケート回答分析システムの全体構成を示すブロック図である。本実施の形態に係るアンケート回答分析システムは、回答者コンピュータ111～11N（Nは正整数）と、ネットワーク12と、データベース13と、テキスト分類エンジン14とから、その主要部が構成されている。

【0014】

回答者コンピュータ111～11Nは、コンピュータ、携帯情報端末、携帯電話機等の各種のメッセージ、メール等の送信機能を持つ機器が該当し、ネットワーク12に接続されている。

【0015】

ネットワーク12には、有線、無線を問わない各種の公衆回線網、専用線網や、LAN（Local Area Network）等の各種のネットワークが含

まれる。

【0016】

データベース13は、ネットワーク12に接続されており、回答者コンピュータ111～11Nからネットワーク12を通じて送信された、複数のアンケート回答者からのアンケート回答文を蓄積する。

【0017】

テキスト分類エンジン14は、データベース13から複数のアンケート回答文を読み出して、アンケート回答文进行分类するルールを導き出し、ルールを要求者に出力する。テキスト分類エンジン14は、データベース13に蓄積されたアンケート回答文の全文に対して形態素解析を行う形態素解析手段15と、テキスト分類エンジン14に対してカテゴリとテキストとを指定させるカテゴリ・テキスト指定手段16と、データベース13から読み込んだ複数のアンケート回答文に対して属性選択を行う属性選択手段17と、属性選択手段17により属性選択された単語を基にテキストとカテゴリとの対応を表すルールを学習するルール学習手段18と、ルールを出力するルール出力手段19とを含んで構成されている。

【0018】

なお、テキスト分類エンジン14は、カテゴリとテキストとの間の対応関係を分類則として学習するエンジンであり、たとえば、Li and Yamanishiにより提案されているエンジン(h. Li and k. Yamanishi: Text Classification Using ESC-based Stochastic Decision Lists, Proceedings of 1999 International Conference on Information & Knowledge Management, pp: 122-130, 1999)を用いることができる。このテキスト分類エンジン14は、基本的に特許第2581196号に開示された「決定リストの生成方法及び装置」の方式を利用したものである。

【0019】

図2は、データベース13に格納されるアンケート回答文の集合の構造を示す。各列はアンケート項目を表し、各行は1人分のアンケート回答文を表す。

【 0 0 2 0 】

図 3 を参照すると、テキスト分類エンジン 1 4 の処理は、形態素解析ステップ 3 1 と、テキストとカテゴリとの指定ステップ 3 2 と、属性選択ステップ 3 3 と、ルール学習ステップ 3 4 と、ルール出力ステップ 3 5 とからなる。

【 0 0 2 1 】

図 4 を参照すると、属性選択ステップ 3 3 のより詳しい処理は、 $\Delta SC(\omega)$ 計算ステップ 4 1 と、属性選択ステップ 4 2 とからなる。

【 0 0 2 2 】

図 5 を参照すると、ルール学習ステップ 3 4 のより詳しい処理は、データ成形ステップ 5 1 と、成長処理ステップ 5 2 と、刈り込み処理ステップ 5 3 とからなる。

【 0 0 2 3 】

図 6 は、テキスト分類エンジン 1 4 による分析結果であるルール形式の知識（確率的決定リスト）の一例を示す図である。

【 0 0 2 4 】

図 7 は、テキスト分類エンジン 1 4 による分析結果であるルール形式の知識（確率的決定リスト）の他の例を示す図である。

【 0 0 2 5 】

次に、このように構成された第 1 の実施の形態に係るアンケート回答分析システムの動作について説明する。

【 0 0 2 6 】

アンケート回答者が、回答者コンピュータ 1 1 1 ~ 1 1 N からアンケート回答文を送信すると、アンケート回答文はネットワーク 1 2 を通じてデータベース 1 3 に蓄積される。ここで、アンケート回答者の数は、N 人とする。この際、アンケート回答文は、自然言語による自由回答記述を含むものであってよい。

【 0 0 2 7 】

テキスト分類エンジン 1 4 は、まず、形態素解析手段 1 5 により、データベース 1 3 に蓄積されたアンケート回答文の全文に対して形態素解析を行う（ステップ 3 1）。

【0028】

次に、テキスト分類エンジン14は、カテゴリ・テキスト指定手段16により、オペレータにアンケート回答文中のカテゴリとテキストとを指定させる（ステップ32）。ここで、カテゴリの指定とは、1つの列の回答に注目して分類するものである。例えば、図2の1列目に注目して、その回答を「A社」と「A社以外」とに分類することが、カテゴリ指定である。また、テキスト指定とは、カテゴリ指定に用いられた列を除いて、分析に用いる列を指定することである。例えば、図2の2列目を選択してテキスト指定を行う。

【0029】

続いて、テキスト分類エンジン14は、属性選択手段17により、データベース13から読み込んだ複数のアンケート回答文に対して属性選択を行う（ステップ33）。ここで、属性選択とは、テキストとカテゴリとの対応を表すのに重要な単語を選択することである。

【0030】

次に、テキスト分類エンジン14は、ルール学習手段18により、属性選択された単語を基にテキストとカテゴリとの対応を表すルールを学習する（ステップ34）。例えば、上記のカテゴリ指定とテキスト指定とをしたときに、図6のようなルールが得られる。

【0031】

図6のルールは、先ず最初の行を読んで、テキストの中で「使いやすい」という単語が現れていれば、そのような回答をした人の92.0%がハイテク企業としてA社を想起している、ということを示す。もし、「使いやすい」という単語が現れていなければ、次に、「未来」と「プライベート」という単語が同時に現れているかどうかをチェックし、現れていれば、そのような回答をした人の87.2%がハイテク企業としてA社を想起している、ということを示す。以下同様に、if-then-else型のルールに従って、上から下へと条件文を読んで行く。このようなルールは、ハイテク企業とハイテク感との間の関係を分かりやすくコンパクトに示すものである。

【0032】

また、別の会社であるB社をとりあげて、「B社」と「B社以外」とにカテゴリ指定したときに、同様な手順で図7のルールが得られたとする。

【0033】

図7のB社のルールを図6のA社のルールと比較すると、A社をハイテク企業と想起する人のハイテク感、使い勝手や嗜好品感覚を重視するのに対し、B社をハイテク企業と想起する人のハイテク感は効率性を重視していることがわかる。このように、ルールの比較により、アンケート回答分析が容易になる。

【0034】

最後に、テキスト分類エンジン14は、ルール出力手段19により、分析結果のルール形式の知識を要求者に出力する（ステップ35）。

【0035】

ここで、ルール形式の知識の例として確率的決定リストをとり上げ、これを生成するための属性選択ステップ33およびルール学習ステップ34についてより詳しく説明する。

【0036】

確率的決定リストとは、図6に示すようなif-then型の確率的ルールの順序付きリストのことである。各確率的ルールは、“ $c = 1 \leftarrow t$ （確率 p ）”の形をとる。ここで、 $c = 1$ は分類の決定、 t は分類の条件、（確率 p ）は確率を表す。

【0037】

まず、属性選択ステップ33について詳しく述べる。

【0038】

属性選択とは、与えられたカテゴリ（例えば、‘A社’と‘A社でない’）に対して、そのカテゴリと深く関係する単語を集めることである。具体的には、図4に示すように、ステップ41で、テキストに出現する単語 ω の各々について、その単語 ω の出現を考慮しない場合のテキスト集合の確率的コンプレキシティ（SCとかく）と考慮する場合のSCとの差 $\Delta SC(\omega)$ を計算し、ステップ42で、その差 $\Delta SC(\omega)$ が、与えられたしきい値 τ より大きければ単語 ω を属性として選択する。

【0039】

以下、SCの具体的な計算の仕方を述べる。入力されたアンケート回答文中のテキストの集合を

$$(d_1, c_1), (d_2, c_2), \dots, (d_m, c_m)$$

と表す。ここで、 d_i は*i*番目のテキストを表し、*i*番目のテキストに現れる単語の列として表現される。 c_i は*i*番目のテキストに対応するカテゴリの値（ラベル）を表し、各 c_i は所与のカテゴリ（‘A社’）に属するならば1、そうでなければ（‘A社でない’）0の値をとる。 m はテキストの数である。

【0040】

さらに、ラベル列を $c^m = c_1, \dots, c_m$ 、テキスト列を $d^m = d_1, \dots, d_m$ と書く。ラベル列 c^m のSCを、数1のように計算する。

【0041】

【数1】

$$SC(c^m) = mH\left(\frac{m^+}{m}\right) + \frac{1}{2} \log \frac{m}{2\pi} + \log \pi$$

【0042】

ここでは、 m^+ はラベル列 c^m 中で値が1であるラベルの数を表す。また、 \log は自然対数を表し、数2とする。

【0043】

【数2】

$$H(z) \stackrel{def}{=} -z \log z - (1-z) \log(1-z)$$

【0044】

たとえば、J. Rissanen, Fisher information and stochastic complexity. IEEE Trans. on Information Theory, 42(1):40-47 (1996) に述べられているように、 $SC(c^m)$ は、ラベル列 c^m を与えられたモ

デル（ここではベルヌイモデル）を用いて記述するのに必要な最短の記述長という意味をもつ。また、

【外 1】

$$c^{m_{\omega}}$$

は、対応するテキスト d_i の中に単語 ω が現れるラベル c_i からなるラベル列であるとする。ここに、 m_{ω} は、

【外 2】

$$c^{m_{\omega}}$$

におけるラベルの数であるとする。すると、

【外 3】

$$c^{m_{\omega}}$$

の SC の値を、数 3 のように計算することができる。

【0 0 4 5】

【数 3】

$$SC(c^{m_{\omega}}) = m_{\omega} H\left(\frac{m_{\omega}^+}{m_{\omega}}\right) + \frac{1}{2} \log \frac{m_{\omega}}{2\pi} + \log \pi$$

【0 0 4 6】

ここでは、 m_{ω}^+ は、

【外 4】

$$c^{m_{\omega}}$$

における値が 1 であるラベルの数を表す。一方、

【外 5】

$$c^{m_{\neg\omega}}$$

は、対応するテキスト d_i に単語 ω が現れないラベル c_i からなるラベル列であるとする。 m_ω は、

【外 6】

$$c^{m_\omega}$$

におけるラベルの数である。

【0047】

【外 7】

$$c^{m_\omega}$$

の SC の値を、数 4 のように計算することができる。

【0048】

【数 4】

$$SC(c^{m_\omega}) = m_\omega H\left(\frac{m_\omega^+}{m_\omega}\right) + \frac{1}{2} \log \frac{m_\omega}{2\pi} + \log \pi$$

【0049】

単語 ω の出現を考慮しない場合の SC と考慮する場合の SC との差 $\Delta SC(\omega)$ は、数 5 のように計算される。

【0050】

【数 5】

$$\begin{aligned} \Delta SC(\omega) &= \frac{1}{m} (SC(c^m) - (SC(c^{m_\omega}) + SC(c^{m_\omega}))) \\ &= \left[H\left(\frac{m^+}{m}\right) - \frac{m_\omega}{m} H\left(\frac{m_\omega^+}{m_\omega}\right) - \frac{m_\omega}{m} H\left(\frac{m_\omega^+}{m_\omega}\right) \right] \\ &\quad - \left[\frac{1}{2m} \log \frac{m_\omega m_\omega \pi}{2m} \right] \end{aligned}$$

【0051】

差 $\Delta SC(\omega)$ の大きい単語 ω は、与えられたカテゴリによく現れる、あるいはほとんど現れない単語である。それらの単語がそのカテゴリと深く関係するとみることができる。そこで、 τ を与えられた閾値とし、 $\Delta SC(\omega) > \tau$ なる単語 ω を属性として選択する。

【0052】

次に、ルール学習ステップ34について詳しく述べる。

【0053】

いま、属性選択された単語が n 個あるとし、 $\omega_1, \dots, \omega_n$ とする。ステップ51で、まず、入力されたテキストの集合を、以下の表現に置き換える。

$(d_1, c_1), (d_2, c_2), \dots, (d_m, c_m)$

【0054】

ここで、各 d_i は、2値ベクトル（一般的に多値の離散ベクトル）

$d_i = (\omega_{i1}, \omega_{i2}, \dots, \omega_{in}) \quad (i = 1, \dots, m)$

を表す。ここで、 ω_{ij} は、属性選択で得られた単語 ω_j が i 番目のテキストに現れれば1、そうでなければ0の値をとる ($j = 1, \dots, n$)。 c_i は i 番目のテキストに対応するカテゴリの値（ラベル）を表し、各 c_i は所与のカテゴリに属するならば1、そうでなければ0の値をとる。 m はテキストの数である。

【0055】

ステップ52では、*if-then-else* 型のルールを選択し、確率的決定リスト A に順番に追加していく。これを「成長」とよぶ。ルールの選択には、例えば、拡張型確率的コンプレキシティ（ ESC とかく）最小原理を用いる。

【0056】

これは、以下のように行う。 k を与えられた正の整数とする。属性選択された単語 ω を基に可能なすべての k 項（単語 ω の k 個までの連言）の集合を T とする。次に、集合 T の項 t の中からテキストに一度も現れないものを取り除く。また、空の確率的決定リスト A を用意する。次に、 ESC の値の減少分がもっとも大きいルールを確率的決定リスト A に順次追加する。

【0057】

ここで、 ESC の計算は、以下のように行う。全データ集合 D を

$(d_1, c_1), (d_2, c_2), \dots, (d_m, c_m)$
 の形のデータの集合とし、ラベル列を $c^m = c_1, \dots, c_m$ とする。ラベル列 c^m のESCの値は、数6のように近似計算することができる。

【0058】

【数6】

$$ESC(c^m) = Loss(c^m) + \lambda \sqrt{m \log m}$$

【0059】

これは、K. Yamanishi, A decision-theoretic extension of stochastic complexity, and its applications to learning, IEEE Trans. Inform. Theory, 44, 1424-1439 (1998) の論文で提案されたオリジナルのESCの1つの近似形である。

【0060】

ここでは、 λ は正の定数を表し、 $Loss(c^m)$ はデフォルトの分類を行う際の誤りの数を表す。デフォルトの分類とは、たとえば、すべてのラベルが0であると仮定することである。 t は集合 T 中の項であるとする。

【0061】

【外8】

c^{m_t}

は、対応するテキスト d_i において項 t が真になるラベル c_i からなるラベル列であるとする。ここに、 m_t は、

【外9】

c^{m_t}

におけるラベルの数であるとする。

【0062】

【外 1 0】

$\text{Loss}(c^{m_t})$

は、項 t による分類を行う際の誤りの数であるとする。一方、

【外 1 1】

$c^{m_{\neg t}}$

は、対応するテキスト d_i において項 t が偽になるラベル c_i からなるラベル列であるとする。ここに、 $m_{\neg t}$ は、

【外 1 2】

$c^{m_{\neg t}}$

におけるラベルの数であるとする。 $\neg t$ は項 t の否定を表す。

【0 0 6 3】

【外 1 3】

$\text{Loss}(c^{m_{\neg t}})$

は、 $\neg t$ による分類を行う際の誤りの数であるとする。

【0 0 6 4】

【外 1 4】

c^{m_t}

と

【外 1 5】

$c^{m_{\neg t}}$

との ESC の値を、それぞれ数 7 および数 8 のように計算することができる。

【0 0 6 5】

【数 7】

$$ESC(c^{m_t}) = Loss(c^{m_t}) + \lambda \sqrt{m_t \log m_t}$$

【0066】

【数 8】

$$ESC(c^{m_{-t}}) = Loss(c^{m_{-t}}) + \lambda \sqrt{m_{-t} \log m_{-t}}$$

【0067】

項 t による分類を行う場合、ESC の値の減少分 $\Delta ESC(t)$ を、数 9 で計算する。

【0068】

【数 9】

$$\begin{aligned} \Delta ESC(t) &= ESC(c^m) - (ESC(c^{m_t}) + ESC(c^{m_{-t}})) \\ &= [Loss(c^m) - Loss(c^{m_t}) - Loss(c^{m_{-t}})] \\ &\quad + [\lambda(\sqrt{m \log m} - \sqrt{m_t \log m_t} - \sqrt{m_{-t} \log m_{-t}})] \end{aligned}$$

【0069】

ESC 最小原理では、 $\Delta ESC(t)$ が最小になるように項 t を選択する。そのような $t = t^*$ が選ばれたならば、これを真とするような全データ集合 D のデータの数を

【外 16】

$$m_t^+$$

とし、そのようなデータのうち、かりに数が多い方のラベルを $c = 1$ として、 $c = 1$ であった数を

【外 17】

$$m_t^{+,+}$$

と、 $c = 0$ であった数を

【外 18】

$$m_{t^*}^-$$

とする。そこで、ルール“ $c = 1 \leftarrow t^*$ (確率 p)”を確率的決定リスト A に追加する。ここで、確率値 p はラプラス推定の方法を用いて、例えば、数 10 により計算する。

【0070】

【数 10】

$$(m_{t^*}^+ + 0.5) / (m_{t^*} + 1)$$

【0071】

集合 T から項 t^* を除いたものを新しい集合 T とし、全データ集合 D から項 t^* を真にする全てのデータを除いたものを新しい全データ集合 D として、全データ集合 D が空になるまで、同じ操作を繰り返す。なお、上記に用いた基準 ESC の代わりに、属性選択に用いた基準 SC を用いることもできる。

【0072】

次に、ステップ 53 では、ステップ 52 で得られた確率的決定リスト A が学習データに過度に適合することがあるので、確率的決定リスト A の最後からルールを 1 つずつ取り除き、 ESC 最小原理の観点からみてこれを取り除かないほうがよいところまで続ける。これを「刈り込み」とよぶ。

【0073】

この場合の ESC 最小原理の適用の仕方を、以下に述べる。まず、ラベル列 c^m の確率的決定リスト A に対する ESC の値を、数 11 のように、確率的決定リスト A におけるすべての項 t に対する ESC の値の和として定義する。

【0074】

【数 1 1】

$$ESC(c^m | A) = \sum_t ESC(c^{m_t})$$

【0 0 7 5】

ただし、

【外 1 9】

ESC(c^{m_t})

は、数 7 のように計算する。次に、ラベル列 c^m と確率的決定リスト A との全 ESC 値を、数 1 2 のように定義する。

【0 0 7 6】

【数 1 2】

$$\begin{aligned} ESC(c^m : A) &= ESC(c^m | A) + \lambda' L(A) \\ &= \sum_t Loss(c^{m_t}) + \lambda \sum_t \sqrt{m_t \log m_t} + \lambda' L(A) \end{aligned}$$

【0 0 7 7】

ここでは、 λ' は正の定数を表す。L(A) は、確率的決定リスト A を符号化する時の符号長である。具体的には、 $L(A) = \log T + \log(T-1) + \dots + \log T(T-i+1)$ と計算する。ただし、T は可能な項 t の数で、i は確率的決定リスト A におけるルールの数である。

【0 0 7 8】

A は刈り込み前の確率的決定リストを、A' は刈り込み後の確率的決定リストを表すものとする。

$$ESC(c^m : A) \leq ESC(c^m : A')$$

【0 0 7 9】

すなわち、

$$ESC(c^m | A') - ESC(c^m | A) \geq \lambda' (L(A) - L(A'))$$

が成り立つ限り、刈り込み処理を続け、この条件が満足されなくなったか、刈り込むべきルールが無くなった時点で得られている確率的決定リストAを出力する。ESCが全体として小さい確率的決定リストAが出力されることになる。

【0080】

第1の実施の形態に係るアンケート回答分析システムでは、ネットワーク12を通じて収集した自然言語による自由回答記述を含むアンケート回答文から自動的に分析結果のルールを導くことができる。

【0081】

第1の実施の形態に係るアンケート回答分析システムにおいて、テキスト分類エンジン14として、例えば、Li and Yamanishiにより提案されているエンジン (h. Li and k. Yamanishi: Text Classification Using ESC-based Stochastic Decision Lists, Proceedings of 1999 International Conference on Information & Knowledge Management, pp: 122-130, 1999) を用いると、 $O(n^k m)$ の計算量で、高速にアンケート回答文からルールを導き出すことができる。ここに、 n はアンケート回答文から属性選択された単語の数、 m はアンケート回答文の数、 k はルール1つの条件にとられる連言に含まれる単語の最大数である。よって、効率的な自動アンケート回答分析が可能となる。また、得られたルールは特定のカテゴリに属するアンケート回答文を `if-then-else` 型のルール形式で、コンパクトにかつ分かりやすく表現している。

【0082】

第1の実施の形態に係るアンケート回答分析システムは、例えば、以下のようなビジネスとして展開できる。企業に関するイメージ調査、特定の商品やサービス等のアンケートを請け負って、図2のようなアンケート項目でアンケートをネットワーク12上で実施し、ネットワーク12を通じてオンラインで集められた自然言語による自由回答記述を含むアンケート回答文をデータベース13に蓄え、そこからアンケート回答文を呼び出して、テキスト分類エンジン14を使用する

ことにより得られるルールを分析結果として要求者に販売する。

【 0 0 8 3 】

(2) 第 2 の実施の形態

図 8 は、本発明の第 2 の実施の形態に係るアンケート回答分析システムの全体構成を示すブロック図である。本実施の形態に係るアンケート回答分析システムは、アンケート回答入力手段 8 1 と、データベース 8 2 と、テキスト分類エンジン 8 3 とから、その主要部が構成されている。

【 0 0 8 4 】

アンケート回答入力手段 8 1 は、ネットワークを介さずにデータベース 8 2 と直接接続されている。

【 0 0 8 5 】

データベース 8 2 は、複数のアンケート回答者からのアンケート回答文を蓄積する。

【 0 0 8 6 】

テキスト分類エンジン 8 3 は、図 1 に示した第 1 の実施の形態に係るアンケート回答分析システムにおけるテキスト分類エンジン 1 4 と全く同様のものである。したがって、対応する部分には同一符号を付して、それらの詳しい説明を省略する。

【 0 0 8 7 】

次に、このように構成された第 2 の実施の形態に係るアンケート回答分析システムの動作について説明する。

【 0 0 8 8 】

アンケート回答入力手段 8 1 は、ネットワークを介さずにデータベース 8 2 と直接接続して、自然言語による自由回答記述を含むアンケート回答文を入力する。

【 0 0 8 9 】

データベース 8 2 は、複数のアンケート回答者からのアンケート回答文を蓄積する。

【 0 0 9 0 】

テキスト分類エンジン 8 3 は、データベース 8 2 から複数のアンケート回答文を読み出して、アンケート回答文を分類するルールを導き出し、分析結果のルールを要求者に出力する。なお、テキスト分類エンジン 8 3 の動作の詳細は、第 1 の実施の形態に係るアンケート回答分析システムにおけるテキスト分類エンジン 1 4 の場合と全く同様であるので、その詳しい説明を割愛する。

【0091】

第 2 の実施の形態に係るアンケート回答分析システムは、例えば、以下のようなビジネスとして展開できる。企業に関するイメージ調査、特定の商品やサービス等のアンケートを請け負って、図 2 のようなアンケート項目でアンケートを実施し、自然言語による自由回答記述を含むアンケート回答文を一度に回収して、これをデータベース 8 2 で蓄え、そこからアンケート回答文を呼び出して、テキスト分類エンジン 8 3 を使用することにより得られる分析結果を要求者に販売する。

【0092】

(3) 第 3 の実施の形態

図 9 は、本発明の第 3 の実施の形態に係るアンケート回答分析システムの全体構成を示すブロック図である。本実施の形態に係るアンケート回答分析システムは、回答者コンピュータ 9 1 1 ～ 9 1 N と、ネットワーク 9 2 と、データベース 9 3 と、テキスト分類エンジン 9 4 と、要求者コンピュータ 9 5 とから、その主要部が構成されている。

【0093】

回答者コンピュータ 9 1 1 ～ 9 1 N は、コンピュータ、携帯情報端末、携帯電話機等の各種のメッセージ、メール等の送信機能を持つ機器が該当し、ネットワーク 9 2 に接続されている。

【0094】

ネットワーク 9 2 には、有線、無線を問わない各種の公衆回線網、専用線網や、LAN 等の各種のネットワークが含まれる。

【0095】

データベース 9 3 は、ネットワーク 9 2 に接続されており、回答者コンピュータ

911～91Nからネットワーク92を通じて送信された、複数のアンケート回答者からのアンケート回答文を蓄積する。

【0096】

テキスト分類エンジン94は、ルール出力手段19が分析結果のルール形式の知識をネットワーク92を通じて送信できる点を除いては、図1に示した第1の実施の形態に係るアンケート回答分析システムにおけるテキスト分類エンジン14と全く同様のものである。したがって、対応する部分には同一符号を付して、それらの詳しい説明を省略する。

【0097】

要求者コンピュータ95は、ネットワーク92を通じてテキスト分類エンジン94に分析結果のルール形式の知識を要求し、テキスト分類エンジン94から分析結果のルール形式の知識をネットワーク92を通じて受信することができる。

【0098】

次に、このように構成された第3の実施の形態に係るアンケート回答分析システムの動作について説明する。

【0099】

アンケート回答者は、回答者コンピュータ911～91Nからネットワーク92を介して自然言語による自由回答記述を含むアンケート回答文を送る。なお、アンケート回答者は、N人とする。

【0100】

データベース93は、ネットワーク92に接続されており、複数のアンケート回答者からのアンケート回答文を蓄積する。

【0101】

テキスト分類エンジン94は、データベース93から複数のアンケート回答文を読み出して、アンケート回答文を分類するルール形式の知識を導き出す。テキスト分類エンジン94は、ネットワーク92に接続し、分析結果のルール形式の知識を要求者コンピュータ95の要求に応じてネットワーク92を通じて配信する。なお、テキスト分類エンジン94の動作の詳細は、分析結果のルール形式の知識がネットワーク92を通じて配信される点を除いて、第1の実施の形態に係る

アンケート回答分析システムにおけるテキスト分類エンジン 1 4 の場合と全く同様であるので、その詳しい説明を割愛する。

【 0 1 0 2 】

第 3 の実施の形態にアンケート回答分析システムは、例えば、以下のようなビジネスとして展開できる。企業に関するイメージ調査、特定の商品やサービス等のアンケートを請け負って、図 2 のようなアンケート項目でアンケートをネットワーク 9 2 上で実施し、ネットワーク 9 2 を通じてオンラインで集められた自然言語による自由回答記述を含むアンケート回答文をデータベース 9 3 で蓄え、そこからアンケート回答文を呼び出して、テキスト分類エンジン 9 4 を使用することにより得られる分析結果を、要求があれば、ネットワーク 9 2 を通じて要求者に配信サービスすることによりビジネスを行う。

【 0 1 0 3 】

(4) 第 4 の実施の形態

図 1 0 は、本発明の第 4 の実施の形態に係るアンケート回答分析システムの全体構成を示すブロック図である。本実施の形態に係るアンケート回答分析システムは、図 1 に示した第 1 の実施の形態に係るアンケート回答分析システムにおいて、データベース 1 3 に接続されたコンピュータ 1 0 1 にテキスト分類エンジンプログラムを記録した記録媒体 1 0 2 を備えるようにしたものである。その他の構成は、第 1 の実施の形態に係るアンケート回答分析システムと全く同様であるので、対応する部分には同一符号を付してそれらの詳しい説明を省略する。

【 0 1 0 4 】

このように構成された第 4 の実施の形態に係るアンケート回答分析システムでは、記録媒体 1 0 2 からコンピュータ 1 0 1 にテキスト分類エンジンプログラムが読み込まれ、形態素解析手段 1 5、カテゴリ・テキスト指定手段 1 6、属性選択手段 1 7、ルール学習手段 1 8、およびルール出力手段 1 9 を含むテキスト分類エンジン 1 4 としてコンピュータ 1 0 1 の動作を制御する。コンピュータ 1 0 1 上でのテキスト分類エンジン 1 4 の動作の詳細は、第 1 の実施の形態に係るアンケート回答分析システムの場合と全く同様になるので、その詳しい説明を割愛する。

【0105】

(5) 第5の実施の形態

図11は、本発明の第5の実施の形態に係るアンケート回答分析システムの全体構成を示すブロック図である。本実施の形態に係るアンケート回答分析システムは、図8に示した第2の実施の形態に係るアンケート回答分析システムにおいて、データベース82に接続されたコンピュータ111にテキスト分類エンジンプログラムを記録した記録媒体112を備えるようにしたものである。その他の構成は、第2の実施の形態に係るアンケート回答分析システムと全く同様であるので、対応する部分には同一符号を付してそれらの詳しい説明を省略する。

【0106】

このように構成された第5の実施の形態に係るアンケート回答分析システムでは、記録媒体112からコンピュータ111にテキスト分類エンジンプログラムが読み込まれ、形態素解析手段15、カテゴリ・テキスト指定手段16、属性選択手段17、ルール学習手段18、およびルール出力手段19を含むテキスト分類エンジン83としてコンピュータ111の動作を制御する。コンピュータ111上でのテキスト分類エンジン83の動作の詳細は、第2の実施の形態に係るアンケート回答分析システムの場合と全く同様になるので、その詳しい説明を割愛する。

【0107】

(6) 第6の実施の形態

図12は、本発明の第6の実施の形態に係るアンケート回答分析システムの全体構成を示すブロック図である。本実施の形態に係るアンケート回答分析システムは、図9に示した第3の実施の形態に係るアンケート回答分析システムにおいて、データベース93に接続されたコンピュータ121にテキスト分類エンジンプログラムを記録した記録媒体122を備えるようにしたものである。その他の構成は、第3の実施の形態に係るアンケート回答分析システムと全く同様であるので、対応する部分には同一符号を付してそれらの詳しい説明を省略する。

【0108】

このように構成された第6の実施の形態に係るアンケート回答分析システムでは

、記録媒体 1 2 2 からコンピュータ 1 2 1 にテキスト分類エンジンプログラムが読み込まれ、形態素解析手段 1 5，カテゴリ・テキスト指定手段 1 6，属性選択手段 1 7，ルール学習手段 1 8，およびルール出力手段 1 9 を含むテキスト分類エンジン 9 4 としてコンピュータ 1 2 1 の動作を制御する。コンピュータ 1 2 1 上でのテキスト分類エンジン 9 4 の動作の詳細は、第 3 の実施の形態に係るアンケート回答分析システムの場合と全く同様になるので、その詳しい説明を割愛する。

【0 1 0 9】

【発明の効果】

【0 1 1 0】

本発明の第 1 の効果は、企業に関するイメージ調査，特定の商品やサービス等のアンケートを請け負って、アンケートをネットワーク上で実施し、ネットワークを通じてオンラインで集められた自然言語による自由回答記述を含むアンケート回答文をデータベースに蓄え、そこからアンケート回答文を呼び出して、テキスト分類エンジンを使用することにより得られるルール形式の知識を分析結果として要求者に販売できることである。

【0 1 1 1】

本発明の第 2 の効果は、企業に関するイメージ調査，特定の商品やサービス等のアンケートを請け負って、アンケートを実施し、自然言語による自由回答記述を含むアンケート回答文を一度に回収して、これをデータベースで蓄え、そこからアンケート回答文を呼び出して、テキスト分類エンジンを使用することにより得られるルール形式の知識を分析結果として要求者に販売できることである。

【0 1 1 2】

本発明の第 3 の効果は、企業に関するイメージ調査，特定の商品やサービス等のアンケートを請け負って、アンケートをネットワーク上で実施し、ネットワークを通じてオンラインで集められた自然言語による自由回答記述を含むアンケート回答文をデータベースで蓄え、そこからアンケート回答文を呼び出して、テキスト分類エンジンを使用することにより得られるルール形式の知識を分析結果として、要求者からの要求に応じてネットワークを通じて配信サービスすることによ

り販売できることである。

【図面の簡単な説明】

【図 1】

本発明の第 1 の実施の形態に係るアンケート回答分析システムの構成を示すブロック図である。

【図 2】

図 1 中のデータベースに蓄積されたアンケート回答文を例示する図である。

【図 3】

図 1 中のテキスト分類エンジンにおける処理を示すフローチャートである。

【図 4】

図 3 中の属性選択ステップのより詳細な処理を示すフローチャートである。

【図 5】

図 3 中のルール学習ステップのより詳細な処理を示すフローチャートである。

【図 6】

図 1 中のテキスト分類エンジンによる分析結果であるルール形式の知識（確率的決定リスト）の一例を示す図である。

【図 7】

図 1 中のテキスト分類エンジンによる分析結果であるルール形式の知識（確率的決定リスト）の他の例を示す図である。

【図 8】

本発明の第 2 の実施の形態に係るアンケート回答分析システムの構成を示すブロック図である。

【図 9】

本発明の第 3 の実施の形態に係るアンケート回答分析システムの構成を示すブロック図である。

【図 1 0】

本発明の第 4 の実施の形態に係るアンケート回答分析システムの構成を示すブロック図である。

【図 1 1】

本発明の第 5 の実施の形態に係るアンケート回答分析システムの構成を示すブロック図である。

【図 1 2】

本発明の第 6 の実施の形態に係るアンケート回答分析システムの構成を示すブロック図である。

【符号の説明】

- 1 2 ネットワーク
- 1 3 データベース
- 1 4 テキスト分類エンジン
- 1 5 形態素解析手段
- 1 6 テキスト・カテゴリ指定手段
- 1 7 属性選択手段
- 1 8 ルール学習手段
- 1 9 ルール出力手段
- 3 1 形態素解析ステップ
- 3 2 テキストおよびカテゴリ指定ステップ
- 3 3 属性選択ステップ
- 3 4 ルール学習ステップ
- 3 5 ルール出力ステップ
- 4 1 $\Delta SC(\omega)$ 計算ステップ
- 4 2 属性選択ステップ
- 5 1 データ成形ステップ
- 5 2 成長処理ステップ
- 5 3 刈り込み処理ステップ
- 8 1 アンケート回答入力手段
- 8 2 データベース
- 8 3 テキスト分類エンジン
- 9 2 ネットワーク
- 9 3 データベース

9 4 テキスト分類エンジン

9 5 要求者コンピュータ

1 0 1, 1 1 1, 1 2 1 コンピュータ

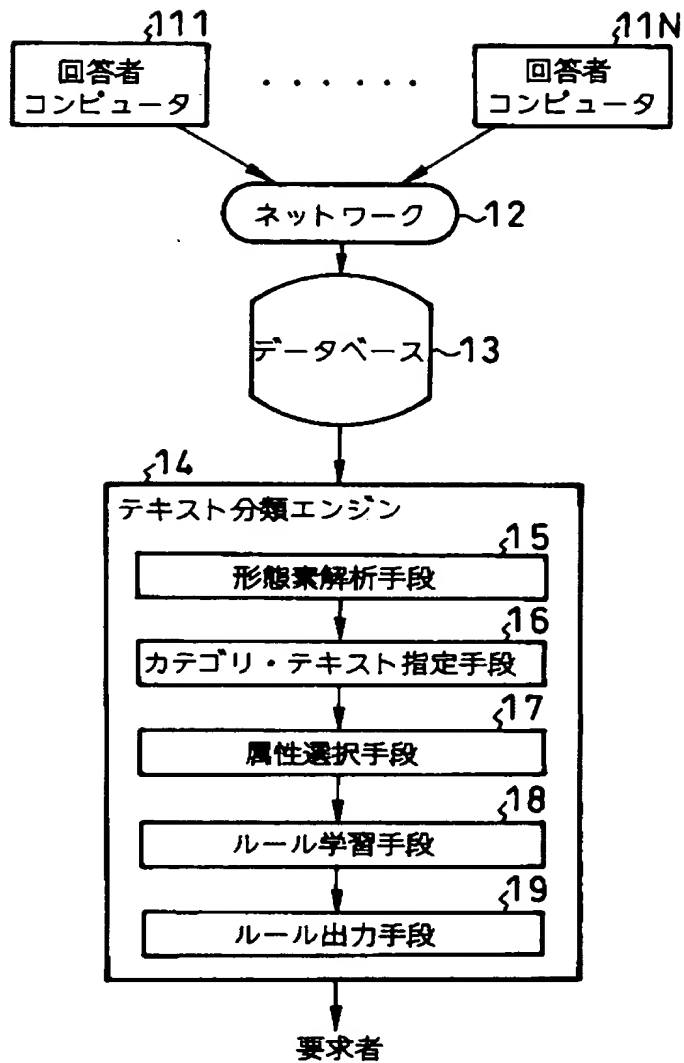
1 0 2, 1 1 2, 1 2 2 記録媒体

1 1 1 ~ 1 1 N 回答者コンピュータ

9 1 1 ~ 9 1 N 回答者コンピュータ

【書類名】 図面

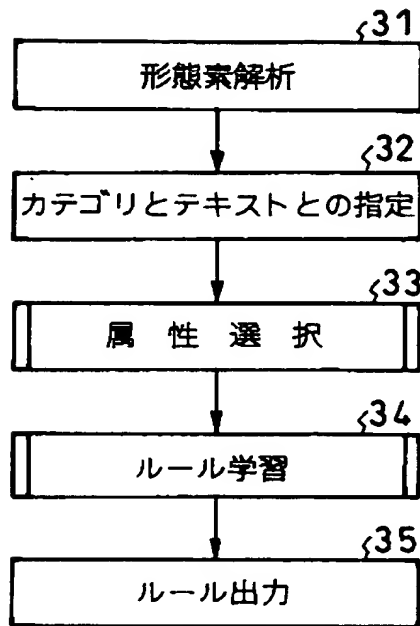
【図 1】



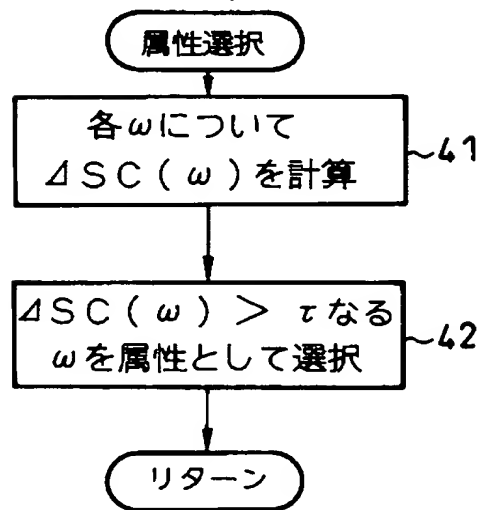
【図 2】

アンケート 回答者	ハイテク企業として どこを想起しますか？	あなたにとってハイテク とは何ですか？	ハイテク商品として 何を思い浮かべますか？
1	A社	進化した未来的な機械	ロボット
2	C社	使いやすく優しいもの	携帯電話
3	A社	高速で高性能な機械	パソコン
...

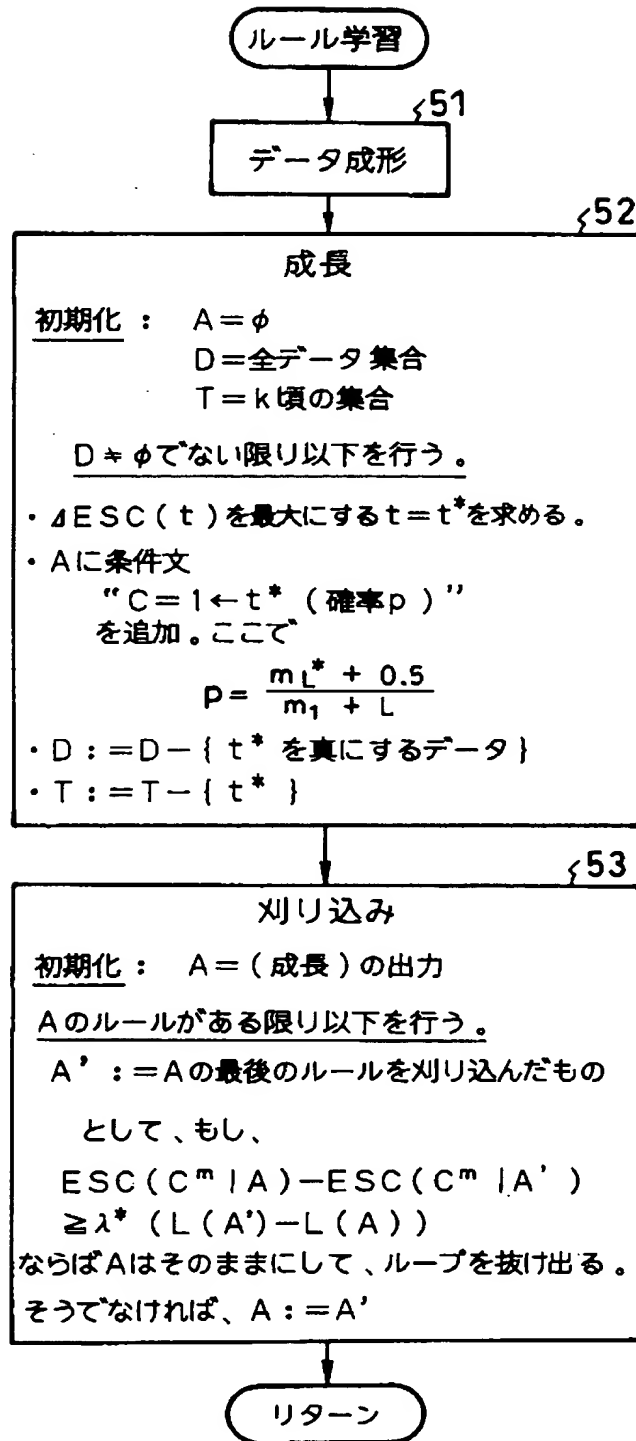
【図 3】



【図 4】



【図 5】



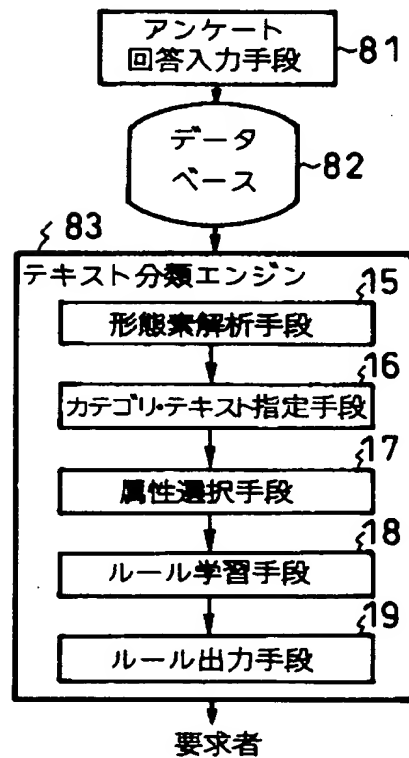
【図 6】

A社 ← 使いやすい〔92.0%〕
 A社 ← 未来&プライベート〔87.2%〕
 A社 ← 疲労&軽減〔78.0%〕
 A社 ← 簡単だ〔65.8%〕
 A社 ← 楽しい〔56.2%〕
 A社でない ← それ以外〔79.4%〕

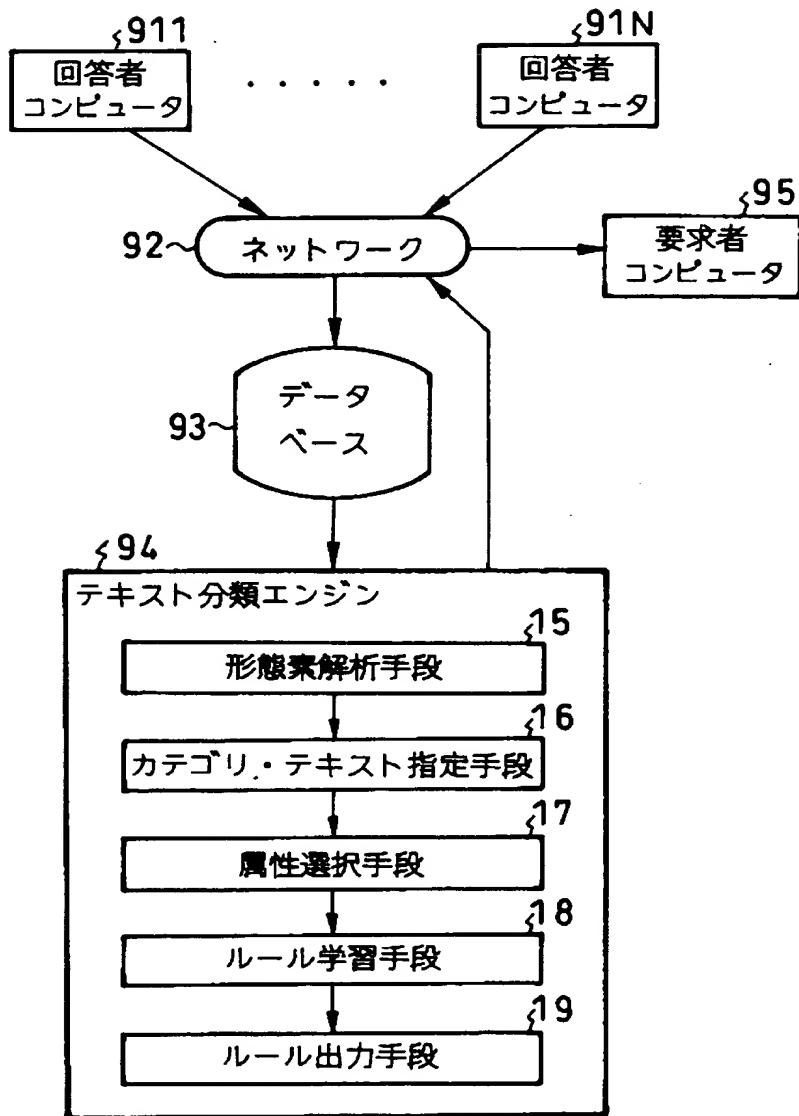
【図 7】

B社 ← 迅速だ〔82.0%〕
 B社 ← 機械&効率〔77.8%〕
 B社 ← 機械&操作〔76.0%〕
 B社 ← 巧妙だ〔60.8%〕
 B社 ← 優れた〔60.2%〕
 B社でない ← それ以外〔76.4%〕

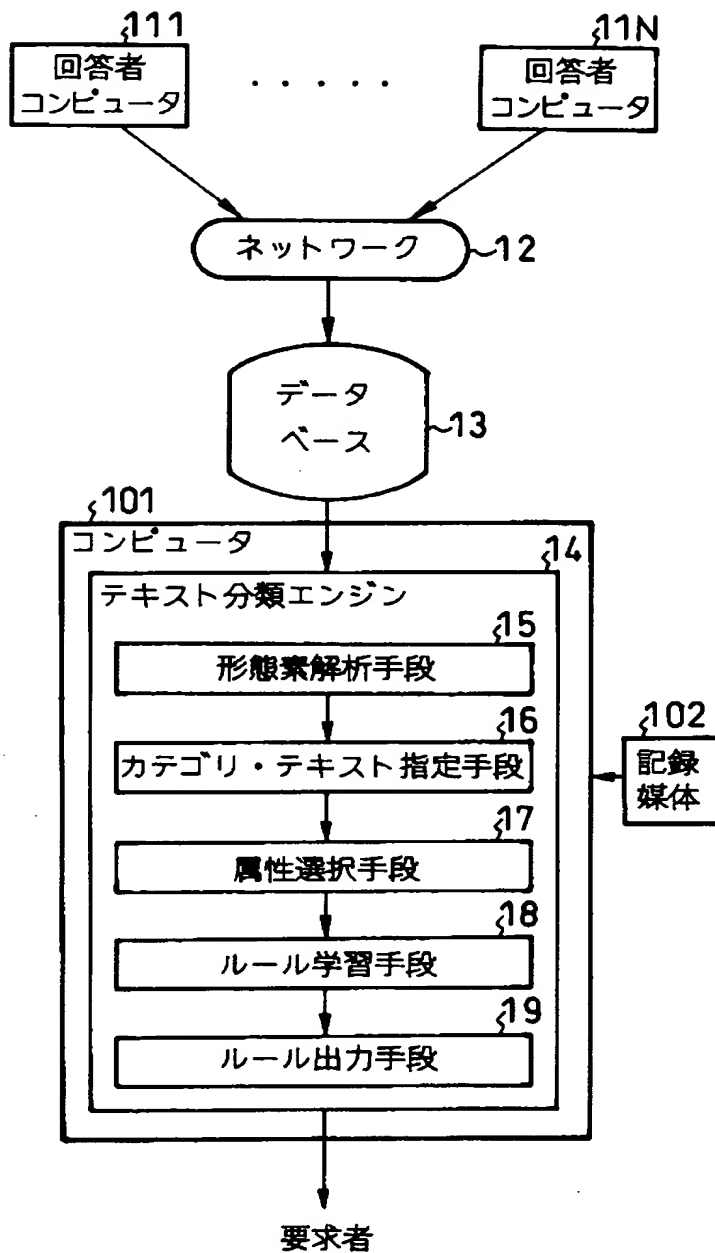
【図 8】



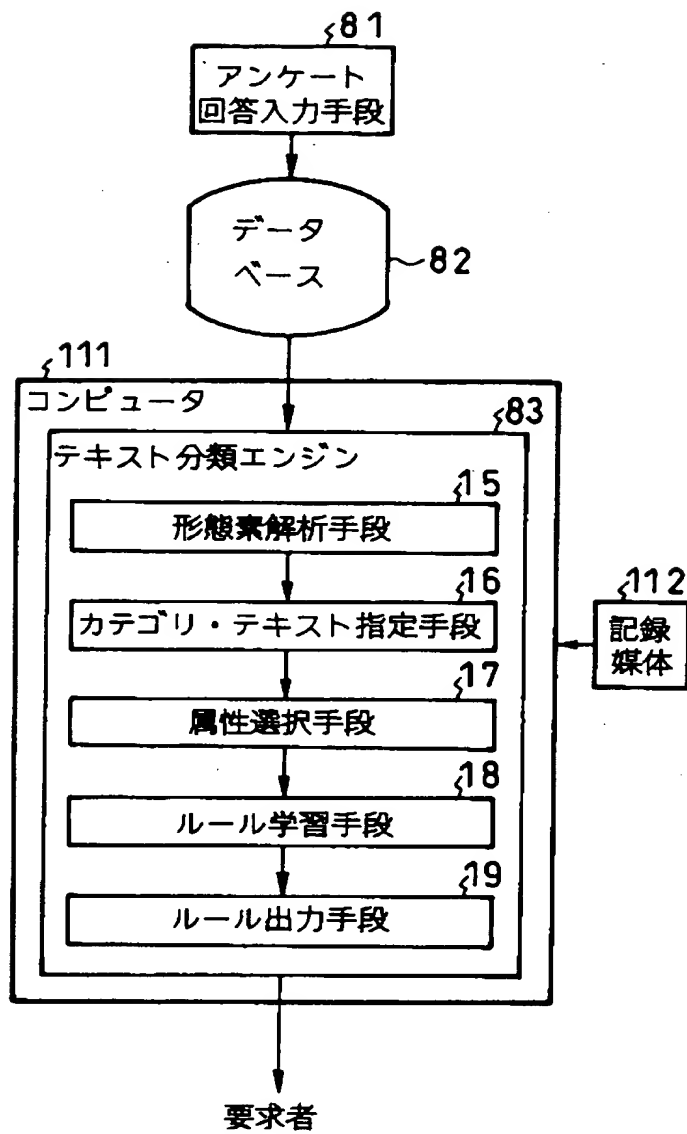
【図9】



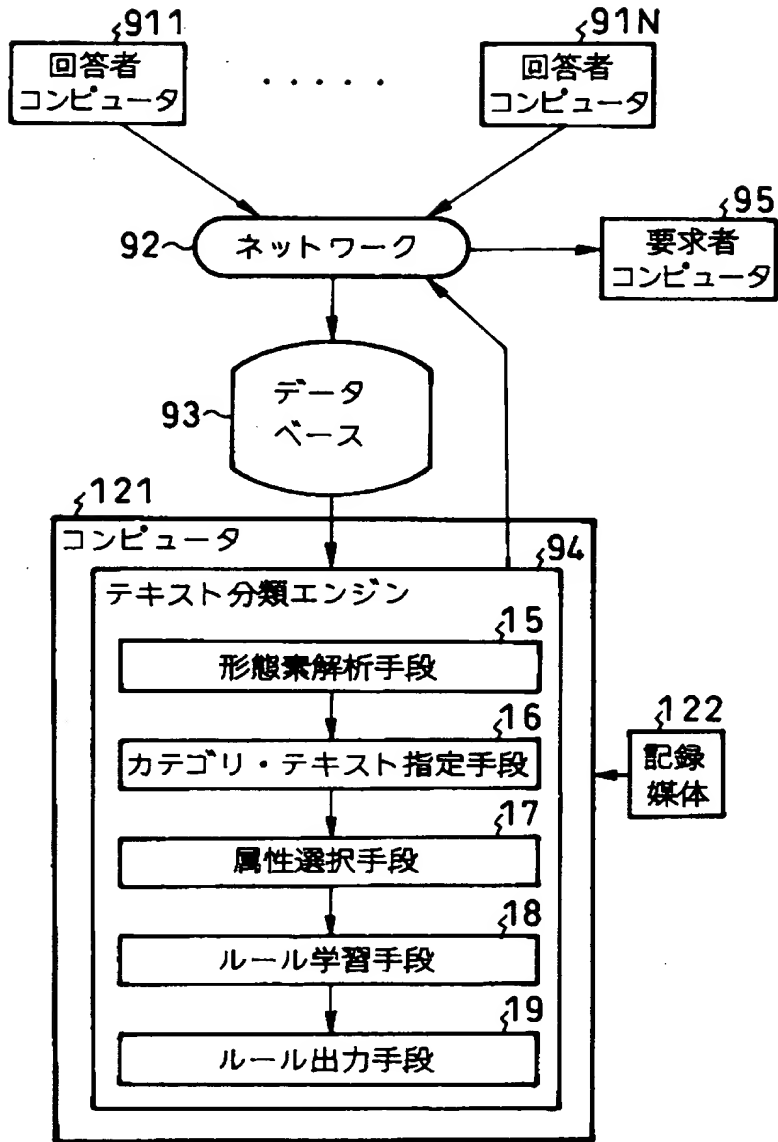
【図10】



【図 11】



【図 12】



【書類名】 要約書

【要約】

【課題】 ネットワークを通じて回収した自然言語による自由回答記述を含むアンケート回答文から、テキスト分類エンジンを用いることにより、アンケート回答分析を自動的に行い、分析結果をルール形式の知識として要求者に出力する。

【解決手段】 アンケート回答者は回答者コンピュータ 1 1 1 ~ 1 1 N からアンケート回答文を送信する。アンケート回答文は、ネットワーク 1 2 を通じてデータベース 1 3 に蓄積される。テキスト分類エンジン 1 4 は、データベース 1 3 から蓄積されたアンケート回答文を読み出して、アンケート回答文を分類するルールを学習して、要求者に出力する。

【選択図】 図 1

認定・付加情報

特許出願の番号	特願2000-071657
受付番号	50000307659
書類名	特許願
担当官	第七担当上席 0096
作成日	平成12年 3月16日

<認定情報・付加情報>

【提出日】	平成12年 3月15日
-------	-------------

出 願 人 履 歴 情 報

識別番号 [000004237]

1. 変更年月日	1990年 8月29日
[変更理由]	新規登録
住 所	東京都港区芝五丁目7番1号
氏 名	日本電気株式会社